FROM GRAPH EMBEDDING TO LKH: BRIDGING LEARNING AND HEURISTICS FOR A STREAMLINED GENERAL TSP SOLVER

Anonymous authors

Paper under double-blind review

Abstract

The Traveling Salesman Problem (TSP) is known as one of the most notorious NP-hard combinatorial optimization problems. In recent decades, researchers from fields such as computer science, operations research, and artificial intelligence including deep learning (DL) have made numerous attempts on the problem. Among the works, the Lin-Kernighan-Helsgaun (LKH) heuristic algorithm is one of the most competent methods for obtaining optimal or near-optimal solutions. Despite the rapid development in DL-based solvers, few of them can defeat LKH in terms of both running efficiency and solution quality across different distributions. In this paper, we would introduce a very novel approach that enhances LKH with graph embedding (GE) techniques in solving general TSP (distances can be non-metric and asymmetric), named as Embed-LKH. It is presented as two stages: i) in the GE stage, it transforms the distances to transition probabilities, then conduct GE given the transition probabilities, and finally it uses the learned embeddings to construct the so-called 'ghost distances'; ii) in the LKH stage, LKH generates candidates based on the ghost distances but searches tours according to the original distances. As the experiments show, compared with the original LKH counterpart, in most cases, our approach can obtain better solutions within the same amount of trials across six distance distributions (non-metric and asymmetric: normal, uniform, exponential, metric and symmetric: Euclidean 2D/10D/50D) and two problem scales (TSP-100/1000). The source files, running scripts, and data are in the anonymous link https://anonymous.4open.science/ r/EmbedLKH-BF80/, which will be made publicly available after the review.

033

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

034 035

036

1 INTRODUCTION

The Travelling Salesman Problem (TSP) is one of the most well-known NP-hard combinatorial optimization (CO) problems. Given a set of nodes and the distances between them, TSP aims to find the shortest tour that visits each node exactly once and returns to the starting node. With applications spanning logistics, telecommunications, and genetics, TSP plays a crucial role in route optimization and network planning (Applegate et al., 2007; Rardin & Rardin, 1998; Matai et al., 2010).

- Isaac Newton.

If I have seen further, it is by standing on the shoulders of giants.

General TSP and its Practical Significance. In this paper, we study the general TSP, where the costs 044 from node to node can be non-Euclidean (are not distances derived from Euclidean coordinates) and asymmetric (the cost from node i to j does not have to be equal to that from j to i). Compared with 046 the symmetric TSP defined by 2D coordinates which is much more commonly studied by literature 047 (Kool et al., 2018; Kwon et al., 2020; Li et al., 2023; Qiu et al., 2022), the general TSP is more 048 ubiquitous in practical applications: For example, the time it takes for vehicles to travel between two places does not entirely depend on the physical distance between the two places, but is highly affected by road congestion and speed limits (65 mph on highway but 25 mph in residential districts). That 051 leads to a non-Euclidean time cost. Meanwhile, the time cost can be also asymmetric due to different congestion situations in different directions (for example, during morning rush hours, the required 052 time of traveling from residential areas to work areas is much higher than the opposite direction), as well as probably different road conditions (some roads are one-way).

054 Research Progress. One algorithm that can-055 not go unmentioned in the context of the TSP 056 is LKH (Lin-Kernighan-Helsgaun (Helsgaun, 057 2000; 2017)), widely regarded as one of the 058 most effective heuristic algorithms for TSP due to its near-optimal solutions, high speed, and good scalability. Huge success has been made 060 in neural solvers for TSP in 2D Euclidean space. 061 Notably, some of the works (Sun & Yang, 2023; 062 Li et al., 2024) achieve a comparable perfor-063 mance with LKH. However, on general TSP, 064 there remains a considerable gap between cur-



Figure 1: Motivation: reducing the difficulty of the input instances helps solvers more easily find better solutions.

rent neural solvers and LKH. In this paper, we are going to bridge this gap by innovatively utilizing graph embedding to enhance LKH, which falls in the so-called Embed-LKH. And for the first time, Embed-LKH outperforms LKH on the general TSP (including the Euclidean ones) within the same amount of trials. Just as we quoted Newton in the preface, we believe that the shortcut to surpassing LKH is to stand on LKH's shoulder, and as shown by experiments, our Embed-LKH does see further.

Motivation. The motivation is straightforward: to reduce the difficulty of input TSP instances. As illustrated in Fig. 1, there are hard cases and easy cases. For hard instances, even strong solvers like LKH struggle to find the optimal solution, whereas for easier instances, a simple greedy algorithm (e.g., nearest neighbor) suffices. Our objective is to transform a difficult instance into an easier one that likely shares the same optimal solution, allowing LKH to find a shorter tour with fewer trials. Importantly, while we adjust the distances during the transformation, the tour length is always computed using the original distances of the hard instance.

We notice that the graph embedding (GE) techniques (e.g., Perozzi et al. (2014); Grover & Leskovec (2016)), originally developed for graph mining tasks (e.g., community discovery, social link prediction, etc.), possess remarkable properties that could be highly effective for general TSP solving, including:
i) Scalability. With a single machine of 128GB memory, GE methods (e.g., Tang et al. (2015); Qiu et al. (2019a)) can easily deal with graphs of millions of nodes. ii) Unsupervised learning. Most GE methods work in an unsupervised style, making it particularly suitable for the general TSP, where optimal solutions are usually inaccessible as supervisory information. iii) Applicability on non-attributed graphs. GE directly learns from graph topology without node attributes. These virtues of GE also motivate us to apply GE in the challenge of solving general TSP.

Our Techniques and Highlighted Features. In our Embed-LKH, we first apply GE techniques on the transition probability matrix transformed from the distance matrix of the general TSP. This allows us to construct a 'ghost distance matrix', a modified distance matrix that we expect will be easier to solve compared to the original distances. We then generate candidates based on the ghost distances while searching tours according to the original distances to ensure the accuracy of the computed tour length. Below, we summarize the main virtues of Embed-LKH:

- Effectiveness. We investigate six distance distributions including three non-metric and asymmetric (normal, uniform, exponential), and three metric and symmetric (Euclidean 2D/10D/50D). Results show that in most cases Embed-LKH can outperform learning-based solvers, as well as the original LKH within the same number of searching trials.
 With a figure result of the same number of searching trials.
- High efficiency and scalability. We utilize closed-form solution for GE, avoiding gradient descent thus achieving higher efficiency of the GE steps. We further propose an asynchonization scheme for Embed-LKH where GE steps and the LKH step are conducted in different threads asynchronously. These technical points endow Embed-LKH with high efficiency and good scalability.
- One-shot. Unlike popular learning-based (Kool et al., 2018; Qiu et al., 2022; Li et al., 2023) and learning-enhanced heuristic (Hudson et al., 2021; Xin et al., 2021) methods for solving TSP, Embed-LKH runs in the one-shot manner, allowing it to be applied to instances of different scales without a typical resource-consuming pre-training stage, while also avoiding generalization issues from training data to testing data.
- Invariance to perturbations. Most supervised (Joshi et al., 2019; Fu et al., 2021) or reinforcement learning-based methods (Kool et al., 2018; Qiu et al., 2022; Kwon et al., 2021) suffer from performance collapse due to the fragility of neural predictions which can be immensely affected by any perturbations on the input distances, largely hindering their applicability to real-world

cases. As we prove in Sec. 6, Embed-LKH is invariant to row-wise distance disturbances and node
 permutations which would not change the optimality of the solutions.

2 RELATED WORKS

For space limit, we put related works in Appendix B where we detailedly analyze the differences between our work and others and summarize the main points in Table 5.

¹¹⁵ 3 BACKGROUNDS AND PRELIMINARIES

We give a formal definition of our studied general TSP in Sec. 3.1. Then in Sec. 3.3, we give preliminaries of the Graph Embedding (GE) techniques. We will elaborate on how to bridge the gap of the two seemingly totally irrelevant worlds in the last.

120 3.1 GENERAL TSP DEFINED BY DISTANCE MATRICES

Definition 1 (General Traveling Salesman Problem). Given a node set \mathcal{V} ($\mathcal{V} = \{1, 2, \dots, N\}$) along with a distance matrix $\mathbf{D}_{i,j} \in [0, +\infty)^{N \times N}$ where the entry $\mathbf{D}_{i,j}$ is the distance from node *i* to *j*, the problem is to find the tour (a Hamiltonian cycle that visits all the nodes exactly once and returns to the starting node) $\tau = (\tau_1, \dots, \tau_N, \tau_1)$ to minimizes the cost $\sum_{i=1}^{N-1} \mathbf{D}_{\tau_i, \tau_{i+1}} + \mathbf{D}_{\tau_N, \tau_1}$. Without losing generality or decreasing the hard level of the problem, we assume that \mathbf{D}_{ij} is a positive value. In this paper, we consider TSP in the general cases, namely

127

147

148 149

153 154

111

128

- Asymmetric. The symmetry, i.e. $\mathbf{D}_{i,j} = \mathbf{D}_{j,i}$, does NOT have to hold.
- Non-Metric. The triangle inequality, i.e. $\mathbf{D}_{i,j} + \mathbf{D}_{j,k} \geq \mathbf{D}_{i,k}$, does NOT have to hold.

129 3.2 LKH ALGORITHM

130 LKH is a highly encapsulated tool tailored for vehicle routing problems including TSP with abundant 131 of technical tricks in it. Diving deep into the working mechanisms of LKH is not the purpose 132 of the paper. Readers who are interested in the technical details of LKH may refer to Appendix C. We highly recommend readers to regard LKH as a black box which works by two steps: 1) 133 LKH GeneCand(D): Given a distance matrix D, it generates *candidates* for each node. Here 134 'candidates' are indicating how likely the edge from the node to a candidate shows in the final 135 solution; 2) LKH_SearchTour(**D**, Candidates): Given distance matrix **D**, it searches tours based 136 on the generated candidates. 137

138 3.3 WORD2VEC-BASED GRAPH EMBEDDING TECHNIQUES

139 Technical Overview of Word2vec. Word2vec (Mikolov et al., 2013) is one of the most well-known 140 algorithms in NLP, winning NeurIPS 2023 Test of Time Awards. It is a shallow neural network model 141 that learns to represent words in a continuous vector space from a given corpus. The training objective 142 of Skip-gram model in word2vec is to predict surrounding *context words* for the *center words*. Word2vec assigns each word i one word embedding $\mathbf{x}_i \in \mathbb{R}^d$, and one hidden embedding $\mathbf{h}_i \in \mathbb{R}^d$. 143 From the natural language text corpus, we can obtain the *empirical co-occurence probability* that a 144 word j shows as the context word of another center word i, denoted by p(j|i), and the frequency that 145 the word i is selected as the center word, denoted by p(i). Then we have the objective: 146

maximize
$$J := \sum_{i,j\in\mathcal{V}} p(i)p(j|i)\log\frac{\exp \mathbf{x}_i^{\top}\mathbf{h}_j}{\sum_{j'\in\mathcal{V}}\exp \mathbf{x}_i^{\top}\mathbf{h}_{j'}},$$
(1)

where \mathcal{V} is the vocabulary (minor notation abuse with node set \mathcal{V} of similar meaning). Notice that the denominator of the objective Eq. 1 is too computationally expensive to be practical. To this end, the negative sampling technique is proposed (Mikolov et al., 2013), whose objective is defined as:

maximize
$$J := \sum_{i,j \in \mathcal{V}} p(i)p(j|i) \Big[\log \sigma(\mathbf{x}_i^{\top} \mathbf{h}_j) + \sum_{s=1}^{S} \mathbb{E}_{j' \sim \mathbb{P}_{neg}} \log \sigma(-\mathbf{x}_i^{\top} \mathbf{h}_{j'}) \Big], \quad (2)$$

where S is the number of negative samples for each positive sample term $\log \sigma(\mathbf{x}_i^{\top} \mathbf{h}_j)$ in the objetive, and \mathbb{P}_{neg} is an empirical distribution of negative samples.

Word2vec-Based GE. Similar to word embedding that aims to represent words as vectors, given an input graph $G = (\mathcal{V}, \mathcal{E})$ (\mathcal{V} : nodes, \mathcal{E} : edges), GE aims to represent the nodes of as vectors for further node-level downstream tasks. Word2vec-based GE methods (e.g., Perozzi et al. (2014); Grover & Leskovec (2016)) first design a random walk strategy (where p(i) and p(j|i) are explicitly





or implicitly defined) to convert the graph-structure data into node sequence data, and then feed the node sequences to word2vec just like dealing with the natural language corpus, after which we would obtain the node embeddings. From this perspective, these word2vec-based GE methods are essentially the art of designing p(i) and p(i|j) for the input graph G. We write such a procedure as the following formula:

$$\mathbf{X} := [\mathbf{x}_i]_{i=1}^N, \ \mathbf{H} := [\mathbf{h}_i]_{i=1}^N \leftarrow \texttt{word2vec}\Big(\mathcal{V}, p(i), p(j|i)\Big)$$
(3)

where $\mathbf{X}, \mathbf{H} \in \mathbb{R}^{d \times N}$ are the node embedding matrix and hidden embedding matrix respectively.

So far, we've already introduced basic techniques of word2vec-based GE. The key step of applying the techniques to general TSP, is to define p(i) and p(j|i) for a given distance matrix **D**. We will explain the methodology in detail in the next section.

186 187

209 210 211

215

174

175

176

177

178

179

181

182

4 EMBED-LKH: GRAPH-EMBEDDING-ENHANCED LKH

188 4.1 METHOD OVERVIEW AND NOTATION EXPLANATIONS

As presented in the backgrounds of last section, GE learns from p(i) and p(j|i). So, step 1 is to transform $\mathbf{D} \in [0, +\infty)^{N \times N}$ in the distance space to a *transition probability* matrix $\mathbf{P} \in [0, 1]^{N \times N}$ in the probability space, whose element \mathbf{P}_{ij} is the transition probability from node *i* to *j*. Step 2 is conducting a word2vec-based GE for **P**. Step 3 is to construct a *ghost distance* matrix $\widehat{\mathbf{D}}$ which is supposed to be an easier case than original distance matrix **D**. And the final step 4 is running LKH to solve **D** but using candidates generated from $\widehat{\mathbf{D}}$.

The method framework is illustrated by Fig. 2, with the complete algorithm presented in Alg. 1 where we give detailed analysis of time complexity of each step. As the supplementary for easy reading, we give the notation list in Table 4, Appendix A.

4.2 STEP 1: TRANSFORM DISTANCES TO TRANSITION PRBABILITY

The Relation between Distances and Probabilities. We consider such a $path^1 \pi = (\pi_1, \pi_2, ..., \pi_M)$ of M nodes. In the distance space, the length of π is defined by the sum of the distances between adjacent nodes in π . While in the probability space, the probability $p(\pi)$ of formulating such a π can be given by cumulative product of the transition probabilities $\mathbf{P}_{\pi_i \pi_{i+1}}$ between adjacent nodes π_i and π_{i+1} in the path. The differences can be mathematically shown by the below Eq. 4:

length (for TSP):
$$L(\pi) = \sum_{i=1}^{M-1} \mathbf{D}_{\pi_i, \pi_{i+1}}$$
, probability (for GE): $p(\pi) = \prod_{i=1}^{M-1} \mathbf{P}_{\pi_i, \pi_{i+1}}$. (4)

Methodology. Based on the observation above, we can define a simple way to transform D to P with a row-wise softmax function over -D as below:

(step 1)
$$\mathbf{P}_{i,j} \leftarrow \text{Softmax}(-\mathbf{D}_{i,:})_j := \frac{\exp(-\mathbf{D}_{i,j})}{\sum_{j'=1}^N \exp(-\mathbf{D}_{i,j'})},$$
 (5)

By Eq. 5, we can see that the smaller $D_{i,j}$ is, the higher transition probability $P_{i,j}$ will be, which is in accordance with the intuition of TSP. More strictly, Eq. 5 also ensures the following Proposition 1:

¹Notation clarification: a tour τ is a special type of paths π where the starting node is also the last node and other nodes show exactly once, see Definition 1.

216 Algorithm 1 Embed-LKH: Graph-Embedding-Enhanced LKH. 217 1: Input: Distance matrix $\mathbf{D} \in \mathbb{R}^{N \times N}$; 218 2: Parameters for GE: embedding dimension d, number of negative samples for each node K. 219 3: Parameters for LKH: the number of runs #runs, the number of maximum trial times #max_trials. 220 4: • Step 1: Distance transformation (Sec. 4.2). Time: $O(N^2)$. 5: $\mathbf{P} \leftarrow diaq(\exp(-\mathbf{D})\mathbf{1})^{-1}\exp(-\mathbf{D})$ 221 ► Transform from distance space to probability space 6: \blacktriangleright Step 2: Graph embedding (Sec. 4.3 and 5.1). 222 7: $\mathbf{Q} \leftarrow f_p(\mathbf{P})$ Time: $O(N^3W)$ for random walk, $O(N^2 \log N)$ for sparsification (optional); 8: $\mathbf{U}_d, \mathbf{\Sigma}_d, \mathbf{V}_d^{\top} \leftarrow \text{SVD}(\log(N\mathbf{Q}/K); d); \mathbf{X} \leftarrow \sqrt{\mathbf{\Sigma}_d}\mathbf{U}_d^{\top}; \mathbf{H} \leftarrow \sqrt{\mathbf{\Sigma}_d}\mathbf{V}_d^{\top}$ **Time:** $O(N^2 \log d + Nd^2)$ 224 9: ► Step 3: Build ghost distance matrix(Sec. 4.4). 225 10: $\hat{\mathbf{P}} \leftarrow \frac{K}{N} \exp(\mathbf{X}^{\top}\mathbf{H}); \hat{\mathbf{P}} \leftarrow \hat{\mathbf{P}} / \operatorname{sum}(\hat{\mathbf{P}}, \dim = 1) \triangleright$ Reconstruct the transition probability by embeddings 226 11: $\widehat{\mathbf{D}} \leftarrow -\log(\widehat{\mathbf{P}}); \widehat{\mathbf{D}} \leftarrow \widehat{\mathbf{D}} - \min(\widehat{\mathbf{D}})$ Build ghost distance matrix D 227 12: ► Step 4: Run LKH with candidates generated from ghost distances (Sec. 4.5) 228 13: Time: $O(N^3)$ for line 14 step and $O(N \text{ #runs} \cdot \text{#max_trials} \cdot \text{#cand})$ for line 15, the same as original LKH 229 14: $Cand \leftarrow LKH_GenCand(\mathbf{D})$ ► Generate candidates from ghost distance matrix 230 15: $\tau^* \leftarrow \text{LKH}_\text{SearchTour}(\mathbf{D}, Cand)$ Search for a best tour τ^* from candidates, 231

Proposition 1 (Distance \mapsto Transitional Probability). We consider a TSP instance attached with a distance matrix **D**. We then define the transition probability **P** by Eq. 5. Observing that the length and probability of a path π defined in Eq. 4 also fit a tour τ , we have:

1. For two tours τ and τ' , if $L(\tau) < L(\tau')$, then $p(\tau) > p(\tau')$.

2. The tour τ of the shortest length $L(\tau)$ is also the tour that has the highest probability $p(\tau)$.

Proof. Substituting π in Eq. 4 with τ , we have

$$p(\tau) = \frac{\exp(-\mathbf{D}_{\tau_N\tau_1})}{\sum_{j'=1}^{N}\exp(-\mathbf{D}_{\tau_Nj'})} \prod_{i=1}^{N-1} \frac{\exp(-\mathbf{D}_{\tau_i\tau_{i+1}})}{\sum_{j'=1}^{N}\exp(-\mathbf{D}_{\tau_ij'})} = \frac{\exp(-L(\tau))}{\prod_{i=1}^{N} \left[\sum_{j'=1}^{N}\exp(-\mathbf{D}_{\tau_ij'})\right]},$$
 (6)

whose denominator is a constant that is independent of τ and only related to **D**, so $p(\tau)$ monotonically decreases with respect to $L(\tau)$. Then the first proposition holds. And the second proposition is obvious with the first proposition proven.

Proposition 1 ensures that the transformation from distance space to probability space by Eq. 5 does not compromise the optimality of the tour, while enabling GE at the same time.

4.3 STEP 2: GRAPH EMBEDDING IN THE TRANSITIONAL PROBABILITY DOMAIN

Recall in Sec. 3.3, we said that word2vec-based GE is the art of designing p(i) and p(j|i). Therefore, the main content of this step is to design a strategy $f_p : \mathbf{P} \mapsto p(i), p(j|i)$, so that we can obtain embeddings **X** and **H** by the below formula following Eq. 1:

(step 2)
$$\mathbf{X}, \mathbf{H} \leftarrow \text{word2vec} (\mathcal{V}, f_p(\mathbf{P})).$$
 (7)

A trivial strategy: p(i) = 1/N, $p(j|i) = \mathbf{P}_{i,j}$ for large instances which saves running time.

Random-walk-based Strategy. Borrowing ideas from deepwalk (Perozzi et al., 2014), we set p(i) = 1/N and $p(j|i) = (\sum_{w=1}^{W} \frac{1}{w} \mathbf{P}^w)_{ij}$, where W is the 'window size', representing how many hops of information will be integrated. The coefficient $\frac{1}{w}$ indicates that the weight of w-hop information reduces as the hop increases. Excessive focus on the one-hop distance between nodes may be detrimental to solving the TSP. Random walk helps to 'see further'.

Sparsification for Large Instances. Before random walk, we select top-K probabilities for each node and mask the rest. It helps to filter out many edges that are unlikely to be candidates in the solution. Note that thought sparsification may improve results significantly, it can also be time-consuming.

After obtaining the embeddings \mathbf{X} and \mathbf{H} , we construct the learnt transition probability matrix $\hat{\mathbf{P}}$ by

(step 3.1)
$$\widehat{\mathbf{P}}_{i,j} \longleftarrow \frac{\exp(\mathbf{x}_i^{\top} \mathbf{h}_j)}{\sum_{j'=1}^{N} \exp(\mathbf{x}_i \mathbf{h}_{j'})},$$
 (8)

232 233

234

235 236

237

238

244

245

246

254 255

256

257

258

259

260

261 262

272

278

279

280

282

283

284

285

287

288 289 290

291

292

293 294

295

296

305 306

312

317

and then construct the '*ghost distance*' matrix $\hat{\mathbf{D}}$ by a similarly reverse procedure of Eq. 5:

(step 3.2)
$$\widehat{\mathbf{D}} \leftarrow -\log(\widehat{\mathbf{P}})$$
, then $\widehat{\mathbf{D}} \leftarrow \widehat{\mathbf{D}} - \min(\widehat{\mathbf{D}})$, (9)

where the second formula aims to set the minimum value of distances as 0. Just as Proposition 1, the procedure of building such a ghost distance matrix from the learnt transitional probability does not hurt the optimality of the tour, as formally presented in Proposition 2. For page limit, the proof is put in Appendix D. Proposition 1 and 2 indicate that the step 1 and 3 will not cause information loss. **Proposition 2**. (Transitional Probability \rightarrow Distance. Reverse to Proposition 1). We consider the

Proposition 2 (Transitional Probability \mapsto Distance, Reverse to Proposition 1). We consider the transition probability matrix $\widehat{\mathbf{P}}$ between a set of nodes \mathcal{V} , with the probability of a path π defined by $\hat{p}(\pi) = \prod_{i=1}^{M-1} \widehat{\mathbf{P}}_{\pi_i,\pi_{i+1}}$ (as Eq. 4). We construct a distance matrix $\widehat{\mathbf{D}}$ by as Eq. 9, then we have:

281 1. For two tours τ and τ' , if $p(\tau) < p(\tau')$, then $\hat{L}(\tau) > \hat{L}(\tau')$ holds in TSP.

2. The tour τ that has the highest probability $\hat{p}(\tau)$ also has the shortest length $\hat{L}(\tau)$ in TSP.

4.5 STEP 4: RUN LKH WITH CANDIDATES GENERATED FROM GHOST DISTANCES

In step 3 we have obtained ghost distances \hat{D} , which we hope is an easier case to solve than original distances D. And in Sec. 3.2 we've introduced how LKH works by two functions LKH_GenCand and LKH_SearchTour. In Embed-LKH, we use the learned ghost distances \hat{D} to generate candidates, then search tours with D, just as blow:

(step 4)
$$Cand \leftarrow \text{LKH}_\text{GenCand}(\widehat{\mathbf{D}}), \quad \tau \leftarrow \text{LKH}_\text{SearchTour}(\mathbf{D}, Cand).$$
 (10)

Since we compute tour length according to **D**, need to accurately calculate the length of the tour to ensure that the output tour is the shortest in the original distance space but not the one in the ghost distance space.

5 TECHNICAL DETAILS FOR HIGHER EFFICIENCY

5.1 CLOSED-FORM SOLUTION OF WORD2VEC BASED ON MATRIX FACTORIZATION

In Sec. 3.3, we've introduced the objective of word2vec (Eq. 1). As proven by previous works (Levy & Goldberg, 2014; Qiu et al., 2018), once p(i) and p(j|i) are explicitly known and no other constraints are introduced during the optimization, the embeddings that are learned from the objective can be optimally obtained via matrix factorization. Compared with the gradient descent-based optimizers, e.g., SGD, the matrix factorization-based optimization is much more efficient.

302 **Methodology.** We first define $y_{ij} := \mathbf{x}_i^\top \mathbf{h}_j$, then we calculate the partial derivative of J with regard 303 to y_{ij} as follows (note that y_ij show not only as the positive sample term, but also as the negative 304 sample term for other positive sample terms):

$$\frac{\partial J}{\partial y_{ij}} = -p(i) \big[p(j|i)\sigma(-y_{ij}) + S\mathbb{P}_{neg}(j)\sigma(y_{ij}) \big].$$
(11)

By setting $\frac{\partial J}{\partial y_{ij}}$ to 0 and using the uniform distribution i.e. $\mathbb{P}_{neg}(j') = 1/N$ for negative sampling (which means all the nodes treated equally as negative samples), we would obtain the optimal solution $y_{ij}^* = \log \frac{p(j|i)}{S\mathbb{P}_{neg}(j)} = \log \frac{Np(j|i)}{S}$. We write p(j|i) in the matrix form as \mathbf{Q} with $\mathbf{Q}_{i,j} = p(j|i)$. Then word2vec is indeed conducting the following matrix factorization from left to right:

$$\log(N\mathbf{Q}/S) \approx \mathbf{X}^{\mathsf{T}}\mathbf{H}.$$
(12)

To obtain the optimal X and H, we can conduct Singular Value Decomposition (SVD) over the LHS term of Eq. 12, and then preserve the highest *d* singular values and corresponding vectors and construct embeddings X and H just as follows:

$$\mathbf{U}_d, \mathbf{\Sigma}_d, \mathbf{V}_d^\top \leftarrow \text{SVD}(\log(N\mathbf{Q}/S); d), \text{ then } \mathbf{X} \leftarrow \sqrt{\mathbf{\Sigma}_d} \mathbf{U}_d^\top \text{ and } \mathbf{H} \leftarrow \sqrt{\mathbf{\Sigma}_d} \mathbf{V}_d^\top.$$
(13)

Randomized SVD. An ordinary procedure of SVD may consume a $O(N^3)$ time. In Embed-LKLH, we have $d \ll N$ (in experiments we set d < 5). In this case, randomized SVD (Halko et al., 2009) with a $O(N^2 \log d + Nd^2)$ time complexity can save much time. Randomized SVD for a given matrix **A** is conducted by two stage: i) Compute an approximate basis for the range of **A**. The goal is to obtain a matrix **B** to make $\mathbf{A} \approx \mathbf{BB}^*$. ii) Use **B** which is much smaller than **A** to compute matrix factorization. Readers who are interested in details are recommended to refer to the official paper.

324 5.2 ASYNCHRONIZE EMBED-LKH

Recall that Embed-LKH has 4 steps, including 3 GE steps (step 1-3) and an LKH step (step 4). To improve running efficiency of Embed-LKH and fully maximize CPU utilization, we asynchronize the GE steps and the LKH step by two separate threads as illustrated in Fig. 3. In this way, when dealing with a batch of tasks, the additional



Figure 3: Asynchronization of Embed-LKH.

computation overheads introduced by the GE steps can be significantly mitigated in terms of their
 impact on the program's runtime efficiency. Go a further step, we can simply increase the number of
 threads as plotted in Fig. 3 to achieve a higher solving speed. In experiments we use 8 threads for the
 LKH step and 8 threads for the GE steps as the default.

337 6 THEORETICAL DISCUSSION

339 6.1 INVARIANCE TO THE INPUT PERTURBATIONS

Abundant works for neural TSP solvers deal with the distance perturbations that would not change
the optimal solutions by data augmentation (Jiang et al., 2022; Kwon et al., 2020; 2021; Ma et al.,
2021). In comparison, Embed-LKH is naturally invariant to some certain perturbations, as presented
by the following Proposition 3 and 4.

Proposition 3 (Row-wise Disturbance Invariance). We define a random row-wise disturbance vector as $\mathbf{b} \in \mathbb{R}^N$, whose entry \mathbf{b}_i represents the disturbance on node *i*'s distances to other nodes. If we disturb the input distance matrix \mathbf{D} by $\mathbf{D}' \leftarrow \mathbf{D} + \mathbf{b}$ ($\mathbf{D}'_{i,j} \leftarrow \mathbf{D}_{i,j} + \mathbf{b}_i$), one can easily prove that such a disturbance would not change the optimality of solutions. We conclude that the GE outputs, embeddings \mathbf{X} and \mathbf{H} , are also invariant to such a disturbance.

Proof. We find that the disturbance works in step 1. For the new distances D', we denote the transition probability as P'. By step 1 (Eq. 5), we have:

351 352 353

354 355

356

349

350

$$\mathbf{P}'_{i,j} = \text{Softmax}(-\mathbf{D}'_{i,:}) = \frac{\exp(-\mathbf{D}_{i,j} - \mathbf{b}_i)}{\sum_{j'=1}^{N} \exp(-\mathbf{D}_{i,j'} - \mathbf{b}_i)} = \frac{\exp(-\mathbf{D}_{i,j})}{\sum_{j'=1}^{N} \exp(-\mathbf{D}_{i,j'})} = \mathbf{P}_{i,j}.$$
 (14)

So, the disturbance would not change the output of step 1, and also has no influence on the final embeddings X and H.

Proposition 4 (Node Permutation Invariance). We define a random permute matrix $\mathbf{M} \in \{0,1\}^{N \times N}$ where each row and each column has exactly one non-zero element. A random permutation over the distance matrix by $\mathbf{D}' \leftarrow \mathbf{M}\mathbf{D}\mathbf{M}^{\top}$ would not change the optimal solution. We conclude that node permutation does not affect the output solution generated by Embed-LKH, either.

³⁶² *Proof.* It is obvious that, by step 1 (Eq. 5), we have the transition probability after permutation ³⁶³ $\mathbf{P}' = \mathbf{MPM}^{\top}$; then by step 2 (Eq. 12 and 13), we have $\mathbf{X}' = \mathbf{MX}$ and $\mathbf{H}' = \mathbf{MH}$; then by step 3 ³⁶⁴ (Eq. 8 and 9), we have $\hat{\mathbf{D}}' = \mathbf{MDM}^{\top}$. Here we find that the ghost distance matrix $\hat{\mathbf{D}}'$ is exactly a ³⁶⁵ randomly permuted $\hat{\mathbf{D}}$. So for the same nodes of different IDs in $\hat{\mathbf{D}}$ and $\hat{\mathbf{D}}'$, the candidates generated ³⁶⁶ by LKH_GenCand would be the same. Also, the function LKH_SearchTour is irrelevant with ³⁶⁷ node orders. So, the solution given by Embed-LKH would not be affected by node permutation. \Box

368 369 370

6.2 EMBED-LKH IS A ONE-SHOT SOLVER

As we presented in Alg. 1, before testing, Embed-LKH is not trained on a training dataset which 371 contains many instances, but conducted right on the testing instance. So, Embed-LKH is a one-372 shot solver. The 'one-shot' property endows the method more other virtues including: i) All-scale 373 applicability. Unlike previous methods that require ensuring that the size of the test data is the same 374 as the size of the instances in the training dataset, our method can be applied to problems of any scale 375 (as long as the machine has enough memory). ii) Low training cost. Unlike previous deep learning 376 methods (e.g., Kwon et al. (2021)), Embed-LKH does not require a significant amount of time for 377 pre-training and has no requirements for hardware (e.g., GPUs). iii) Generalizability. Embed-LKH does not have the issue of generalization from training data to testing data.

³⁷⁸ 7 EXPERIMENTS

380 7.1 EXPERIMENTAL SETUP

Hardware. Experiments for neural methods are conducted on a single NVIDIA RTX3090 24GB
 GPU with AMD 3970X 32-Core CPU. Our Embed-LKH and other heuristics (LKH, EAX, etc.) are
 conducted on a machine with Intel Xeon W-3175X CPU and 128GB memory.

384 Baselines. Neural solver: MatNet (Kwon et al., 2021) is the only neural solver runnable for general 385 TSP to our best knowledge. Heuristic solver: LKH (Helsgaun, 2017) is a very strong heuristic for 386 its high efficiency, good scalability, and impressive performance. Learning-enhanced heuristic 387 solver: VSR-LKH (Zheng et al., 2021) improves LKH's heuristic strategy with RL-driven policies. 388 The improved versions (VSR-LKH-V2/3) do not support general (asymmetric) TSP solving as we 389 have tried. We **DO NOT** compare with the exact solver **Gurobi** since it runs out of memory for TSP instances of 100 nodes or more on our machine; and we **DO NOT** compare with **Concorde** since it 390 does not support distance matrices as the input. Some newly proposed methods are not open-sourced 391 yet (e.g., (Drakulic et al., 2024)) or do not provide an ATSP solver (e.g., (Drakulic et al., 2023; Ye 392 et al., 2024b)), so we do not compare with them. 393

Default Parameters. For fairness, for LKH, VSR-LKH, and Embed-LKH, we set runs=1 (number of running iterations for each instance), max_candidates=6 (maximum number of candidates for each node). Other detailed settings (e.g., the search algorithm, the way to compute candidates, etc.) are set as the default of the original LKH program. Embed-LKH runs with 16 threads (8 for GE steps and 8 for the LKH step) as the default. On TSP-100, we set random walk window size W = 3 without sparsification; on TSP-100 we set W = 1 and adopt sparsification with K = 100.

400 **Testing Data Generation.** Before introducing data generation we should emphasize again that 401 Embed-LKH is a one-shot solver so the training data is also the testing data. We consider the scales fo 100 nodes and 1000 nodes. The investigated distance distributions include three non-euclidean 402 asymmetric types of distances (normal, uniform, exponential) and three types of Euclidean (metric 403 and symmetric) distances: i) Euclidean 2D/10D/50D. First we generate 2D/10D/50D Euclidean 404 random coordinates in $[0, 1]^{2/10/50}$. Then we compute the pair-wise distances, and finally re-scale 405 them to the range of [0, 1e4]. ii) Normal. We Generate distance $D_{i,j}$ from normal distribution 406 $\mathcal{N}(0;1)$, then re-scale them to the range of [0, 1e4]. iii) Exponential. We generate distance $\mathbf{D}_{i,i}$ 407 from exponential distribution ($\lambda = 0.5$), then re-scale them to the range of [0, 1e6]. iv) Uniform. We 408 generate distances from uniform distribution in [0, 1e5]. All distances are rounded to integers (since 409 the float number would cause inaccurate tour lengths). 410

Metrics. Average length (abbr. 'Length'). We report the average length of the found tours of all the instances. Optimal Gap (abbr. 'Opt. Gap'). We take LKH with the largest number of trials of each scale as the reference to compute the performance gap of all compared methods. Time. We report the running time of the solver. For Embed-LKH, it includes the time of all the 4 steps.

415 7.2 MAIN RESULTS

We give results of TSP-100 (100nodes, 1000 instances) in Table 1 and TSP-1000 (1000nodes, 100 instances) in Table 2. We set d = 3 for TSP-100 and d = 1 for TSP-1000. Embed-LKH outperforms VSR-LKH and MatNet consistently, and We provide analysis from the following aspects:

Observation 1: In most cases, Embed-LKH can find better solutions than LKH within the same
 number of trials. It demonstrates the effectiveness of Embed-LKH's generating candidates from the
 ghost distances instead of the original distances.

Observation 2: Embed-LKH is better at non-metric and asymmetric data, and also competent
on high-dimensional (dimension ≥ 10) Euclidean distances. The gap between other methods and
LKH on the normal, uniform, and exponential distances is significantly higher than on the Euclidean
distances. We also observe that Embed-LKH is at the top level on high-dimensional Euclidean data.
The results demonstrate the significance of GE.

Observation 3: LKH is still competent on Euclidean data. As shown in the results, on TSP-100, when max_trials \leq 1000, LKH is better than Embed LKH on Euclidean 2D data; and LKH is also superior to Embed-LKH on TSP-1000 2D Euclidean data with max_trials=10. This demonstrates the effectiveness of original LKH in 2D Euclidean problems. Theoretically, this is because the search algorithm of LKH (λ -opt) by exchanging edges in a tour has a strong physical meaning in its

Satting	Mathod		Normal			Uniform		E	xponential	
Setting	wiethou	Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓
	MatNet	204521.857	4.753%	2m15s	793294.879	381.074%	1m45s	2323096.288	1244.204%	6m3s
may trials	LKH	196221.741	0.502%	1m18s	176221.947	6.865%	1m24s	186287.126	7.791%	3m18s
=10	VSR-LKH Embed-LKH	197116.059 195629.278	0.960% 0.198%	1m25s 21s	182244.461 170434.696	10.518% 3.356%	1m28s 20s	192711.099 179324.534	11.508% 3.762 %	2m1s 41s
may trials	LKH	195369.814	0.065%	1m25s	166124.554	0.742%	1m27s	174487.398	0.963%	2m9s
=100	VSR-LKH Embed-LKH	199476.359 195311.846	2.169% 0.036%	3m48s 22s	196072.837 165745.077	18.904% 0.512%	2m22s 23s	208767.031 173894.059	20.798% 0.620%	2m48s 43s
may trials	LKH	195274.407	0.017%	3m17s	165180.832	0.170%	3m7s	173104.776	0.163%	6m2s
=1000	VSR-LKH Embed-LKH	201260.070 195249.138	3.082% 0.004%	21m19s 40s	206323.724 165018.977	25.120% 0.072%	12m2s 40s	219028.828 172987.579	26.736% 0.095%	11m52s 1m
may trials	LKH	195241.982	-	17m39s	164900.749	-	16m18s	172823.147	_	16m38s
=10000	VSR-LKH Embed-LKH	195234.123	-0.004%	2m	Out o	of Time (time> -0.053%	·3h) 3m9s	172727.728	-0.055%	3m18s
		E	uclidean 2D		Eu	uclidean 10D		Eu	clidean 50D	
		Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓
	MatNet	93240.803	53.309%	1m43s	489627.558	37.171%	1m43s	664242.653	3.724%	5m57s
max trials	LKH	61586.028	1.261%	1m16s	359495.360	0.714%	1m29s	641949.717	0.243%	1m31s
=10	VSR-LKH Embed-LKH	62861.306 61598.954	3.358% 1.282%	1m34s 40s	361055.618 358744.417	1.151% 0.503%	1m53s 24s	642914.388 641480.050	0.394% 0.170%	2m14s 22s
	LKH	60903.262	0.139%	1m50s	357832.551	0.248%	1m59s	640998.944	0.095%	2m14s
=100	VSR-LKH Embed-LKH	63106.889 60910.506	3.762% 0.151%	3m19s 34s	362271.816 357564.439	1.492% 0.173%	3m28s 26s	643562.573 640753.293	0.495% 0.056%	5m11s 25s
	LKH	60823.766	0.008%	7m45s	357235.218	0.081%	7m4s	640581.210	0.030%	8m23s
=1000	VSR-LKH Embed-LKH	63163.195 60822.693	3.854% 0.006%	56m46s 1m20s	363138.487 357076.636	1.734% 0.036%	25m30s 1m5s	644148.481 640450.215	0.587% 0.009%	35m21s 1m
may trials	LKH	60818.954	-	58m53s	356947.480	-	46m7s	640392.055	-	44m43s
=10000	VSR-LKH Embed-LKH	60818.497	-0.001%	7m41s	Out o 356880.873	of Time (time> -0.019%	·3h) 6m49s	640320.954	-0.011%	6m51s

Table 1: Results on TSP-100. Bold: methods yielding best results with the same setting.

Table 2: Results on TSP-1000, 100 instances.

Catting	Mathad	Normal		Uniform			Exponential			
Setting	Method	Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓
max_trials =10	LKH VSR-LKH Embed-LKH	1826493.81 1829168.70 1809316.69	_ 0.146% -0.940%	6m10s 6m27s 42m2s	250579.77 245652.75 206660.97	-1.966% - 17.527%	7m31s 6m23s 42m13s	180838.91 181457.10 174166.86		6m16s 6m4s 8m10s
		Eu	clidean 2D		Ει	clidean 10D		Euc	clidean 50D	
		Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓	Length↓	Opt. Gap↓	Time↓

design for Euclidean 2D data: on a 2D Euclidean plane, the sum of the lengths of the diagonals of a quadrilateral can never be shorter than the sum of the lengths of its opposite sides. We also observe that on Euclidean 10D/50D, the gap between LKH and Embed-LKH is significantly smaller than the gap on the non-metric data, indicating LKH's advantages on the metric Euclidean data.

Observation 4: Asynchronized Embed-LKH is efficient in most cases. As shown by Table 1, on TSP-100, the running time of asynchronized Embed-LKH is less than most baselines. While on TSP-1000, time consumed by the additional steps of GE becomes non-ignorable. We leave the discussion of how much asynchronization improves the efficiency in different cases in Sec. 7.3.

7.3 PARAMETER ANALYSIS AND ABLATION STUDIES

Window size W. Experiments are conducted on TSP-1000 of the exponential distribution with a varying W. Results are plotted in Fig. 5 a), showing that the best window size is about W = 3. When W = 1, it means that no random walk is conducted. W = 3 is better than W = 1 indicates that a proper setting of random walk is beneficial.

Embedding dimension d. We conduct experiments on TSP-100, exponential distribution, with a varying d. We plot results in Fig. 5 b), showing that either a too high or too low d can cause a decline in performance. d = 3 is the best. This also implies that simply augmenting the transition probability matrix through random walks is not enough. The dimensionality reduction step in GE is necessary.

Number of negative samples S. Experiments are conducted on TSP-100 of the exponential distri-bution with a varying S. Results plotted in Fig. 5 c) show that under this certain setting, S = 20



Figure 5: Ablation studies over parameters d, W, S, and the sparsification operation. All metrics are better when lower.

might be a good choice. However, from the absolute variation of the average length, the influence of parameter S on the results is not significant when $S \ge 5$.

499 Studies over the sparsification technique (in 500 step 2). We conduct study from the following 501 two aspects: i) Effectiveness on different data. 502 We compare Embed-LKH with and without spar-503 sification on TSP-1000 of different distributions. 504 We plot the percentage change in length before 505 and after applying the sparsification technique in Fig. 5 d). The results show that sparsification 506



has a positive effect on some data distributions (uniform, normal, Euclidean 10D and 50D) but a 507 negative effect on others (exponential, Euclidean 2D). This implies that in some scenarios, sparsi-508 fication might incorrectly mask edges that have the potential to be part of the optimal solution. ii) 509 **Parameter analysis.** We run on TSP-1000 uniform data with different sparsification factor K, with 510 results plotted in Fig. 4. Results show that when K < 200, the average tour length is considerably 511 smaller than running without sparsification, demonstrating the effectiveness of sparsification in some 512 cases and also highlighting the importance of selecting a proper K. iii) Efficiency-Effectiveness 513 trade-off. The sparsification step may lead to significant time consumption. In our experiments 514 of Fig. 4, we found that K = 100 would consumes 2533 seconds, which is more than 5 times of 515 running without sparsification (461 seconds).

516 Efficiency improvement by multi-thread

asynchronization. In Sec. 5.2 we have introduced the asynchronization of GE steps and LKH step. We run Embed-LKH with and without asynchronization and report the running time as shown in Table 3. In the table, 'Thread=2' means we use one

Data and Setting	No Asynchorinaztion	Threads=2	Threads=16
TSP-100, Normal max_trials=100	192s	172s	26s
TSP-1000, Uniform max_trials=10	20574s (24x of Thread=16)	16576s (19x of Thread=16)	857s

Table 3: Running time w/ and w/o multi-thread asynchronization

thread for GE steps and one for LKH; 'Thread=16' means 8 for GE and 8 for LKH. We see that when
 the number of threads increases, the runtime correspondingly decreases. We also observe that the
 acceleration brought by asynchronization is more significant on larger-scale problems (TSP-1000).

Other findings. We provide some interesting findings on the differences between ghost distances and original distances in Appendix E, which align well with our motivation shown in Fig. 1.

528 529 8 CONCLUSIONS

In this paper, we present a graph-embedding-enhanced LKH (Embed-LKH) algorithm for general
 TSP solving. This is the first technology to utilize graph embedding techniques to assist in solving
 the general TSP, and also the first learning-enhanced one that outperforms LKH on the problem of
 general TSP showing in the ML community, to our best knowledge.

Limitations and Future Works. i) The setting of parameter d and the strategy p(j|i) (especially the inside parameter, window size W) in the GE steps may require careful consideration. We leave automated parameter selection as the future work. ii) There is still a significant room for improving the efficiency of Embed-LKH by more tightly coupling the two functions LKH_GenCand and LKH_SearchTour.

539

526

527

495

540 **ETHICS STATEMENT** 541

542 The method proposed in this paper enhances the performance of general TSP solving via the combi-543 nation of neural graph embedding technique and the heuristic algorithm LKH. A dataset comprising 544 synthetic instances of TSPs with different distributions will be released upon publication. To our best 545 knowledge, no potential harmful impacts can be observed that need be otherwise stated.

Reproducibility Statement

The hardware, parameters, settings, and the generation of datasets, are provided in Sec. 7.1. An anonymous repository is presented at the end of the abstract as to demonstrate basic implementation and results for reviewing. Source code and datasets shall be made public upon publication.

REFERENCES

546 547

548 549

550

552 553

554

565

570

571

572

577

- 555 David L. Applegate, Robert E. Bixby, Vašek Chvátal, and William J. Cook. 2007. 556
- Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with global structural 557 information. In CIKM, pp. 891-900. ACM, 2015. 558
- 559 Yuxiao Dong, Nitesh V. Chawla, and Ananthram Swami. metapath2vec: Scalable representation learning for 560 heterogeneous networks. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2017. 561
- 562 Darko Drakulic, Sofia Michel, Florian Mai, Arnaud Sors, and Jean-Marc Andreoli. Bq-nco: Bisimulation 563 quotienting for efficient neural combinatorial optimization, 2023. URL https://arxiv.org/abs/ 564 2301.03313.
- Darko Drakulic, Sofia Michel, and Jean-Marc Andreoli. Goal: A generalist combinatorial optimization agent 566 learner, 2024. URL https://arxiv.org/abs/2406.15079. 567
- Xingbo Du, Junchi Yan, Rui Zhang, and Hongyuan Zha. Cross-network skip-gram embedding for joint network 569 alignment and link prediction. TKDE, 2022.
 - Zhang-Hua Fu, Kai-Bin Qiu, and Hongyuan Zha. Generalize a small pre-trained model to arbitrarily large tsp instances. In Proceedings of the AAAI conference on artificial intelligence, pp. 7474–7482, 2021.
- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In Proceedings of the 22nd 573 ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 855–864. ACM, 2016. 574
- 575 Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Stochastic 576 algorithms for constructing approximate matrix decompositions. arXiv preprint arXiv:0909.4061, 909, 2009.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In Advances 578 in neural information processing systems, pp. 1024–1034, 2017. 579
- 580 Keld Helsgaun. An effective implementation of the lin-kernighan traveling salesman heuristic. European journal of operational research, 126(1):106-130, 2000.
- 582 Keld Helsgaun. An extension of the lin-kernighan-helsgaun tsp solver for constrained traveling salesman and 583 vehicle routing problems. Roskilde: Roskilde University, 12, 2017. 584
- Benjamin Hudson, Qingbiao Li, Matthew Malencia, and Amanda Prorok. Graph neural network guided local 585 search for the traveling salesperson problem. arXiv preprint arXiv:2110.05291, 2021. 586
- 587 Yuan Jiang, Yaoxin Wu, Zhiguang Cao, and Jie Zhang. Learning to solve routing problems via distributionally 588 robust optimization. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pp. 9786-9794, 2022. 589
- 590 Chaitanya K Joshi, Thomas Laurent, and Xavier Bresson. An efficient graph convolutional network technique 591 for the travelling salesman problem. arXiv preprint arXiv:1906.01227, 2019. 592
- Minsu Kim, Junyoung Park, and Jinkyoo Park. Sym-nco: Leveraging symmetricity for neural combinatorial 593 optimization. Advances in Neural Information Processing Systems, 35:1936–1949, 2022.

594 595	Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In <i>Inter-</i> national Conference on Learning Representations, 2017. URL https://openreview.net/forum?
596	id=SJU4ayYg1.
597 598	Wouter Kool, Herke Van Hoof, and Max Welling. Attention, learn to solve routing problems! <i>arXiv preprint</i>
599	<i>urxiv</i> .1803.08475, 2018.
600	Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Youngjune Gwon, and Seungjai Min. Pomo: Policy
601	optimization with multiple optima for reinforcement learning. Advances in Neural Information Processing
602	<i>Systems</i> , 33:21188–21198, 2020.
603	Yeong-Dae Kwon, Jinho Choo, Ilioo Yoon, Minah Park, Duwon Park, and Youngiune Gwon, Matrix encoding
604 605	networks for neural combinatorial optimization. Advances in Neural Information Processing Systems, 34: 5138–5149 2021
606	5156 5119, 2021.
607	Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. Advances in neural information processing systems, 27, 2014.
608	
609 610	Yang Li, Jinpei Guo, Runzhong Wang, and Junchi Yan. From distribution learning in training to gradient search in testing for combinatorial optimization. In <i>Thirty-seventh Conference on Neural Information Processing</i> <i>Systems</i> , 2023. URL https://openreview.net/forum?id=JtF0ugNMv2.
011	
012	Yang Li, Jinpei Guo, Runzhong Wang, and Junchi Yan. From distribution learning in training to gradient search
613	in testing for combinatorial optimization. Advances in Neural Information Processing Systems, 36, 2024.
614	Weivi Liu, Pin-Yu Chen, Sailung Yeung, Toyotaro Suzumura, and Lingli Chen, Principled multilaver network
615	embedding. In 2017 IEEE International Conference on Data Mining Workshops (ICDMW), pp. 134–141.
616	IEEE, 2017.
617	Oiang Ma Suwan Ga Danyang Ha Darshan Thakar, and Iddo Drori. Combinatorial optimization by graph
618	nointer networks and hierarchical reinforcement learning arXiv preprint arXiv:1911.04936. 2019
619	
620	Yining Ma, Jingwen Li, Zhiguang Cao, Wen Song, Le Zhang, Zhenghua Chen, and Jing Tang. Learning to
621 622	iteratively solve routing problems with dual-aspect collaborative transformer. Advances in Neural Information Processing Systems, 34:11096–11107, 2021.
623	Raiesh Matai Surva Prakash Singh and Murari Lal Mittal Traveling salesman problem: an overview of
624	applications, formulations, and solution approaches. <i>Traveling salesman problem, theory and applications</i> , 1
625	(1):1–25, 2010.
626	
627	and phrases and their compositionality. In NIPS 2013
628	and phrases and their compositionarity. In Wir 5, 2015.
629 630	Yimeng Min, Yiwei Bai, and Carla P. Gomes. Unsupervised learning for solving the travelling salesman problem, 2024. URL https://arxiv.org/abs/2303.10538.
631	
632	Yuichi Nagata and Shigenobu Kobayashi. A powerful genetic algorithm using edge assembly crossover for the
633	travening salesman problem. <i>INFORMS Journal on Computing</i> , 25(2):540–505, 2015.
634	Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. Asymmetric transitivity preserving graph
635	embedding. In Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and
636	<i>data mining</i> , pp. 1105–1114, 2016.
637	Bryan Perozzi Rami Al-Rfou and Steven Skiena, Deenwalk, Online learning of social representations. In
638	Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining,
639	pp. 701–710. ACM, 2014.
640	Fighers On Varias Dang Has Ma Fight L Varian Ware and F. T. M. (1991)
641	factorization: Unifying deepwalk line and node2vec. In WSDM 2018
642	actorization. Onlying deep wark, nile, pie, and node2vee. In <i>w5Dim</i> , 2010.
643	Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Chi Wang, Kuansan Wang, and Jie Tang. Netsmf: Large-scale
644	network embedding as sparse matrix factorization. In <i>The World Wide Web Conference</i> , pp. 1509–1520. ACM, 2010
645	2019a.
646	Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Chi Wang, Kuansan Wang, and Jie Tang. Netsmf: Large-scale
647	network embedding as sparse matrix factorization. In <i>The World Wide Web Conference</i> , pp. 1509–1520. ACM, 2019b.

648 649	Ruizhong Qiu, Zhiqing Sun, and Yiming Yang. DIMES: A differentiable meta solver for combinatorial optimization problems. In <i>Advances in Neural Information Processing Systems</i> 35, 2022.
650 651 652	Meng Qu, Jian Tang, Jingbo Shang, Xiang Ren, Ming Zhang, and Jiawei Han. An attention-based collaboration framework for multi-view network representation learning. In <i>CIKM</i> , 2017.
653 654	Ronald L Rardin and Ronald L Rardin. <i>Optimization in operations research</i> , volume 166. Prentice Hall Upper Saddle River, NJ, 1998.
655 656 657	Leonardo FR Ribeiro, Pedro HP Saverese, and Daniel R Figueiredo. struc2vec: Learning node representations from structural identity. In <i>Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , pp. 385–394. ACM, 2017.
658 659 660	Martin JA Schuetz, J Kyle Brubaker, and Helmut G Katzgraber. Combinatorial optimization with physics-inspired graph neural networks. <i>Nature Machine Intelligence</i> , 4(4):367–377, 2022.
661 662 663 664	Jingyan Sui, Shizhe Ding, Boyang Xia, Ruizhi Liu, and Dongbo Bu. Neuralgls: learning to guide local search with graph convolutional network for the traveling salesman problem. <i>Neural Comput. Appl.</i> , 36 (17):9687–9706, October 2023. ISSN 0941-0643. doi: 10.1007/s00521-023-09042-6. URL https://doi.org/10.1007/s00521-023-09042-6.
665 666	Zhiqing Sun and Yiming Yang. DIFUSCO: Graph-based diffusion solvers for combinatorial optimization. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023.
667 668 669	Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In <i>Proceedings of the 24th International Conference on World Wide Web</i> , pp. 1067–1077. International World Wide Web Conferences Steering Committee, 2015.
670 671 672	Anton Tsitsulin, Davide Mottin, Panagiotis Karras, and Emmanuel Müller. Verse: Versatile graph embeddings from similarity measures. In <i>Proceedings of the 2018 World Wide Web Conference</i> , pp. 539–548. International World Wide Web Conferences Steering Committee, 2018.
673 674 675 676	Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/29921001f2f04bd3baee84a12e98098f-Paper.pdf.
677 678 679	Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. In <i>KDD</i> , pp. 1225–1234. ACM, 2016.
680 681 682	Liang Xin, Wen Song, Zhiguang Cao, and Jie Zhang. Neurolkh: Combining deep learning model with lin- kernighan-helsgaun heuristic for solving the traveling salesman problem. In <i>Advances in Neural Information</i> <i>Processing Systems</i> , volume 34, 2021.
683 684 685	Hao Xiong, Junchi Yan, and Li Pan. Contrastive multi-view multiplex network embedding with applications to robust network alignment. In <i>Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery; Data Mining</i> , 2021.
686 687	Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? <i>arXiv</i> preprint arXiv:1810.00826, 2018.
688 689 690	Haoran Ye, Jiarui Wang, Zhiguang Cao, Helan Liang, and Yong Li. Deepaco: neural-enhanced ant systems for combinatorial optimization. <i>Advances in Neural Information Processing Systems</i> , 36, 2024a.
691 692 693	Haoran Ye, Jiarui Wang, Helan Liang, Zhiguang Cao, Yong Li, and Fanzhang Li. Glop: Learning global partition and local construction for solving large-scale routing problems in real-time. In <i>Proceedings of the AAAI</i> <i>Conference on Artificial Intelligence</i> , 2024b.
694 695 696	Ziwei Zhang, Peng Cui, Xiao Wang, Jian Pei, Xuanrong Yao, and Wenwu Zhu. Arbitrary-order proximity preserved network embedding. In <i>Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining</i> , pp. 2778–2786, 2018.
697 698 699 700	Jiongzhi Zheng, Kun He, Jianrong Zhou, Yan Jin, and Chu-Min Li. Combining reinforcement learning with lin-kernighan-helsgaun algorithm for the traveling salesman problem. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 35(14):12445–12452, May 2021. doi: 10.1609/aaai.v35i14.17476. URL https://ojs.aaai.org/index.php/AAAI/article/view/17476.

702 703		Appendix	
704 705	C	ONTENTS	
706 707			
708	A	Notations	15
709 710	В	Related Works	15
711		B.1 TSP Solvers	15
712		P 2 Word2vac Based Graph Embedding	16
713 714		B.2 word2vec-Based Graph Enibedding	10
715	С	Some Technical Key Points of LKH	17
717 718	D	Proof to Proposition 2	17
719	I.	Other Latenation E's dimen	17
720	E	Other Interesting Findings	1/
721			
722			
723			
724			
725			
720			
728			
729			
730			
731			
732			
733			
734			
735			
730			
738			
739			
740			
741			
742			
743			
744			
745			
740			
748			
749			
750			
751			
752			
753			
754			
755			

Table 4: Main notations and descriptions.

Notations	Descriptions
N	The number of nodes in TSP
$\mathbf{D} \in \mathbb{R}^{N imes N}$	Distance matrix
$\widehat{\mathbf{D}} \in \mathbb{R}^{N imes N}$	Ghost distance matrix
L(au)	Length of a tour τ or a path π in TSP
$\mathbf{P} \in \mathbb{R}^{N imes N}$	Transition probability transformed from D
p(j i)	Empirical co-occurance probability that node/word j
	shows in the neighborhood of node/word i
Q	p(j i) in the matrix form
$\widehat{\mathbf{P}} \in \mathbb{R}^{N imes N}, \widehat{p}(j i)$	Transition probability learned by embeddings, in matrix
	form and scalar form respectively
$p(\pi)$	Probability that a path π (or a tour τ) forms in the proba-
	bility space
J	Objective for GE
$\mathbf{x}_i, \mathbf{h}_i \in \mathbb{R}^d$	Node embedding and hidden embedding of node i
$\mathbf{X},\mathbf{H}\in\mathbb{R}^{d imes N}$	Node embedding and hidden embedding in the matrix form
y_{ij}	The value of $\mathbf{x}_i^{T} \mathbf{h}_j$
$\mathbf{U}_d, \mathbf{V}_d \in \mathbb{R}^{N imes d}, \mathbf{\Sigma}_d \in \mathbb{R}^{d imes d}$	Results of SVD

A NOTATIONS

We list the main notations in Table 4 for reference in reading.

B RELATED WORKS

B.1 TSP SOLVERS

787 Exact Solvers. Solvers for linear programming and mixed integer linear programming, e.g., Gurobi
 788 and CPLEX, can be used to solve general TSP with exact optimal solutions as output. These methods
 789 can be time-consuming in real-world applications where the scale of instances may go large. Concorde
 790 is a solver developed for TSP but it fails to deal with general TSP where only distances are available.

Heuristic Solvers. Some trivial heuristic algorithms include nearest neighbor (NN), furthest insertion
 (FI), etc. After decades of optimizations, strong heuristics such as Lin-Kernighan-Helsgaun (LKH by
 Helsgaun (2017)) and EAX (Nagata & Kobayashi, 2013)) have been well developed and are able to
 give satisfying results with good efficiency and scalability.

Neural Solvers. Compared with the heuristics and exact solvers, neural solvers are believed to be more efficient when the solving procedure is run parallelly on GPUs. Also, the learning-based approaches are believed to have a stronger capability in solving TSP instances from a specific known distribution. Most neural TSP solvers are developed for TSP with coordinates (Kool et al., 2018; Kwon et al., 2020; Li et al., 2023; Qiu et al., 2022; Sun & Yang, 2023; Vinyals et al., 2015; Min et al., 2024). This is due to the current lack of neural networks that can effectively learn in scenarios with **only** pairwise feature information, such as the distances in TSP. MatNet (Kwon et al., 2021) is one of the works that can be applied to BiTSP. It proposes a Transformer-based solver for asymmetric TSP (ATSP) whose input is a distance matrix. However, the heavy-encoder heavy-decoder neural architecture and the DRL-based training put limitations to its scalability and consequently its applicability to real-world problems whose scales are varying. Recently, several works (Drakulic et al., 2023; 2024; Ye et al., 2024b) conducted experiments on ATSP and showed promising results. Yet, such part regarding general TSP solving have not open-sourced.

808 Developing neural solvers that combine LKH and neural networks is also an important line of 809 research. VSR-LKH (Zheng et al., 2021) improves LKH's heuristic strategy with RL-driven policies, which achieves performance improvement on 2D Euclidean TSP compared with vanilla LKH. It is a

000	-	• Embed-LKH	 General 	GE+LKH, • one-shot	✓ (Upon publication)
835		VSR-LKH (Zheng et al., 2021)	• General	RL+LKH, • one-shot	✓ /
834	neuristic solver	NeuroLKH (Xin et al., 2021)		SGN+LKH, pre-trained	✓
833	Learning-enhanced	DeepACO (Ye et al., 2024a)		GNN+ACO, pretrained	\checkmark
832		NeuralGLS (Sui et al., 2023)	Euclidean	GNN+GLS, pretrained	×
031		GNNGLS (Hudson et al., 2021)		GNN+GLS, pretrained	✓
001		MatNet (Kwon et al., 2021)		Transformer+RL, pre-trained	\checkmark
830		GLOP (Ye et al., 2024b)	 General 	Divide-Conquer+RL, pre-trained	No ATSP solver proposed
829		BQ-NCO (Drakulic et al., 2023)		GCN+Transformer+RL, pre-trained	ATSP solver not provided
828		QUBO (Schuetz et al., 2022)		GNN+UL, pre-trained	✓
827		UTSP (Min et al., 2024)		SAG+UL, pre-trained	\checkmark
826		DIFUSCO(Sun & Yang, 2023) T2T(Li et al., 2024)		Diffusion+Generative, pre-trained	\checkmark
o∠4 825	Neural solver	GCN(Joshi et al., 2019) Att-GCN(Fu et al., 2021)		GNN+SL, pre-trained	\checkmark
020		DIMES (Qiu et al., 2022)	Euclidean	GNN+meta RL, pre-trained	
823		SYM-NCO (Kim et al., 2022)			
021		AM (Kool et al., 2018) POMO (Kwon et al., 2020)		GAT+DRL pre-trained	✓ ✓
821		Ma et al. (2019)		Si TUTKE, pre trained	
820		PointerNet (Vinyals et al., 2015)		GPN+RI pre-trained	\checkmark
910	Heuristic solver	EAX(Nagata & Kobayashi, 2013)	Symmetric	Genetic algorithm	\checkmark
818		LKH (Helsgaun, 2017)	General	See Sec. C	\checkmark
817		Concorde	Euclidean		
816	Exact solver	CPLEX		IVA	not open-sourced
815	Exact solver	Gurobi	General	N/A	Available but
814					
813	Solver Type	Method	Distance Type	Description	Open-sourced
812					

Table 5: Technical comparison of Embed-LKH and other TSP Solvers. It should be noted that although many newly proposed neural solvers for TSP has not yet been open-sourced.

learning-based solver yet without neural networks. NeuroLKH (Xin et al., 2021), improves LKH with the Sparse Graph Network (SGN) aimed at generating better edge candidates as a substitution for LKH's α -nearest measure. However, SGN relies on node coordinates as input, making NeuroLKH inapplicable to general TSP.

We provide a technical comparison between Embed-LKH and all other TSP solvers in Table 5,
demonstrating the significance of our work.

845 846

847

848

837 838 839

840

841

842

B.2 WORD2VEC-BASED GRAPH EMBEDDING

849 We mainly discuss about Graph Embedding (GE) methods based on the language model word2vec (?). In this branch, GE methods have been designed for different types of networks, including homoge-850 neous networks (Perozzi et al., 2014; Grover & Leskovec, 2016; Tang et al., 2015), heterogeneous 851 networks (Dong et al., 2017) and multiplex networks (Liu et al., 2017; Qu et al., 2017; Xiong 852 et al., 2021). By designing the input node sequences, GE methods achieve to preserve specific 853 types of information, e.g., low-order proximity (LINE (Tang et al., 2015),) structure similarity 854 (struc2vec (Ribeiro et al., 2017)), versatile similarity measures (VERSE (Tsitsulin et al., 2018)), 855 and cross-network alignment relations (CENALP (Du et al., 2022)). Another line of GE research is 856 developing methods based on matrix factorization. GraRep (Cao et al., 2015) solves the embedding 857 by matrix factorization for random walk and Skip-Gram while also taking high-order proximity 858 information into consideration. NetMF (Qiu et al., 2018) proposes to unify some word2vec-based 859 GE methods including LINE, DeepWalk, and node2vec, etc. within a matrix factorization framework, 860 and NetSMF (Qiu et al., 2019b) treats the network embedding task as sparse matrix factorization. 861 AROPE and HOPE (Zhang et al., 2018; Ou et al., 2016) propose to preserve multi-order proximity by matrix factorization with some mathematical tricks. Methods based on deep neural networks (Wang 862 et al., 2016) especially graph neural networks (Hamilton et al., 2017; Kipf & Welling, 2017; Xu et al., 863 2018) are also influential in literature of GE.

864 C SOME TECHNICAL KEY POINTS OF LKH

For readers not interested in how LKH works, just skip this part and regard LKH as a black box that works efficiently to find near-optimal tours for TSP. In LKH, there are three main components working in series that contribute to its strong performance,

870 Step 1. Compute the α -Nearest Measure. Suppose L(T) is the length of the minimum 1-tree (Def. 2) of graph $G = (\mathcal{V}, \mathcal{E})$ and $L(T^+(i, j))$ is the length of the minimum 1-tree containing edge 872 (i, j), the α -measure of edge (i, j) can be calculated as $\alpha(i, j) = L(T^+(i, j)) - L(T)$. α -measures 873 are used for specifying the edge candidate set which consists of p edges with the smallest α -measures 874 connected to each node (p = 5 as default).

B75 Definition 2 (1-tree). Let $G = (\mathcal{V}, \mathcal{E})$ be a undirected weighted graph where $\mathcal{V} = \{1, 2, \dots, n\}$ is the set of nodes and $\mathcal{E} = \{(i, j) | i \in \mathcal{V}, j \in \mathcal{V}\}$ is the set of edges. A 1-tree for a graph $G = (\mathcal{V}, \mathcal{E})$ is a spanning tree on the node set $\mathcal{V} \setminus \{1\}$ combined with two edges from \mathcal{E} incident to node 1. A minimum 1-tree T is a 1-tree of minimum length.

Step 2. Node Penalties. LKH uses a subgradient optimization technique to obtain a penalty π_i over each node *i*, then modify the distance as $C_{ij} \leftarrow C_{ij} + \pi_i + \pi_j$. This operation changes the distances but would not affect the optimal solution.

Step 3. Search Solutions by λ **-Opt.** λ edges in the current tour will be exchanged by another set of λ edges, all of which should be included in the candidate sets of the original λ edges, to improve the tour until no such exchanges can be found.

D PROOF TO PROPOSITION 2

Here we provide the proof to Proposition 2. By Eq. 4 and Eq. 9, we have

$$\hat{L}(\tau) = \sum_{i=1}^{M-1} \widehat{\mathbf{D}}_{\tau_i, \tau_{i+1}} + \widehat{\mathbf{D}}_{\tau_M, \tau_1} - N \min(\widehat{\mathbf{D}})$$

$$= -\log\left(\widehat{\mathbf{P}}_{\tau_M, \tau_1} \prod_{i=1}^{M-1} \widehat{\mathbf{P}}_{\tau_i, \tau_{i+1}}\right) - N \min(\widehat{\mathbf{D}})$$

$$= -\log \widehat{p}(\tau) - N \min(\widehat{\mathbf{D}})$$
(15)

by which it is obvious that $\hat{p}(\tau)$ monotonically decreases with respect to $\hat{L}(\tau)$. Then the two propositions are obviously right.

E OTHER INTERESTING FINDINGS

One of the questions that we are curious about is that what are the differences between the ghost distances and the original distances. Our answer is that **GE make the non-metric data more like metric data.** We investigate the TSP-100 instances of the normal distribution, we find that the number of non-metric cases $(\mathbf{D}_{i,j} + \mathbf{D}_{j,k} < \mathbf{D}_{i,k})$ decreases significantly in the ghost distance matrix after GE (from more than 10000 cases to less than 1000 cases). As the experiments show, LKH is strong at solving metric data. So, to LKH, making the non-metric data more like metric data may be beneficial to non-metric TSP solving. And that is exactly our motivation as shown in Fig. 1.